


HW2 Assignment



- Home Page
- Recent Search Engine Articles
- Schedule of Lectures
- Assignments
- Special Resources
- Course Grading
- Course Materials

CS572 Course Assignments

Last Modified: December 15, 2017

Homework 1: Comparing Search Engine Results

- [Search Engine Comparison Exercise](#)
- [Grading Guidelines](#)
- Homework #1 Due Jan. 29th

Homework 2: Web Crawling

- [\[Instructions for Installing Eclipse and crawler4j\]](#)
- [\[Flowchart for Crawler4j\]](#)
- [\[Web Crawler Exercise\]](#)
- [\[Grading Guidelines\]](#)
- Homework #2 Due Feb. 21st

- Involves
 1. Java programming
 - I assume all of you know how to program in Java!
 2. Eclipse Software Development Environment
 3. crawler4j, an open source java web crawler
 4. a crawl and analysis of a web site and an analysis of the crawl

What is Eclipse?

- Eclipse started as a proprietary IBM product (IBM Visual age for Smalltalk/Java)
 - Embracing the open source model IBM opened the product up
- Open Source
 - It is a general purpose open platform that facilitates and encourages the development of third party plug-ins
- Best known as an Integrated Development Environment (IDE)
 - Provides tools for coding, building, running and debugging applications
- Originally designed for Java, now supports many other languages
 - Good support for C, C++
 - Python, PHP, Ruby, etc...

Prerequisites for Running Eclipse

- Eclipse is written in Java and will thus need an installed JRE (Java Runtime Environment) or JDK (Java Development Kit) in which to execute
 - JDK recommended

Obtaining Eclipse

- Eclipse can be downloaded from...
 - <http://www.eclipse.org/downloads/packages/>
- Eclipse comes bundled as a zip file (Windows) or a tarball (all other operating systems)

Installing Eclipse

- Simply unwrap the zip file to some directory where you want to store the executables

- The document

“Instructions for Installing Eclipse and crawler4j”

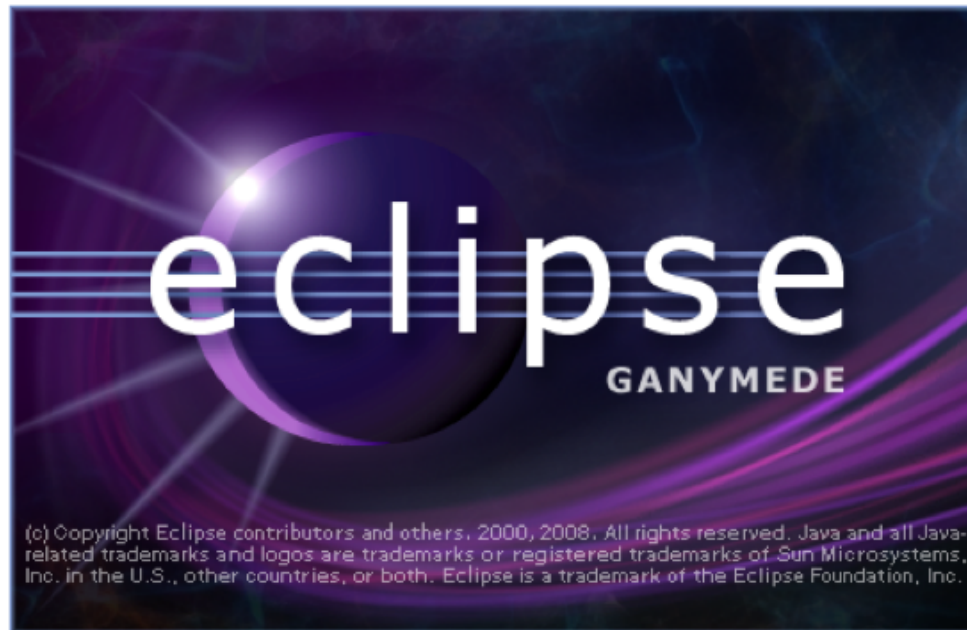
- located at

<http://www-scf.usc.edu/~csci572/2018Spring/hw2/Crawler4jinstallation.pdf>

describes the installation for Windows and Macs

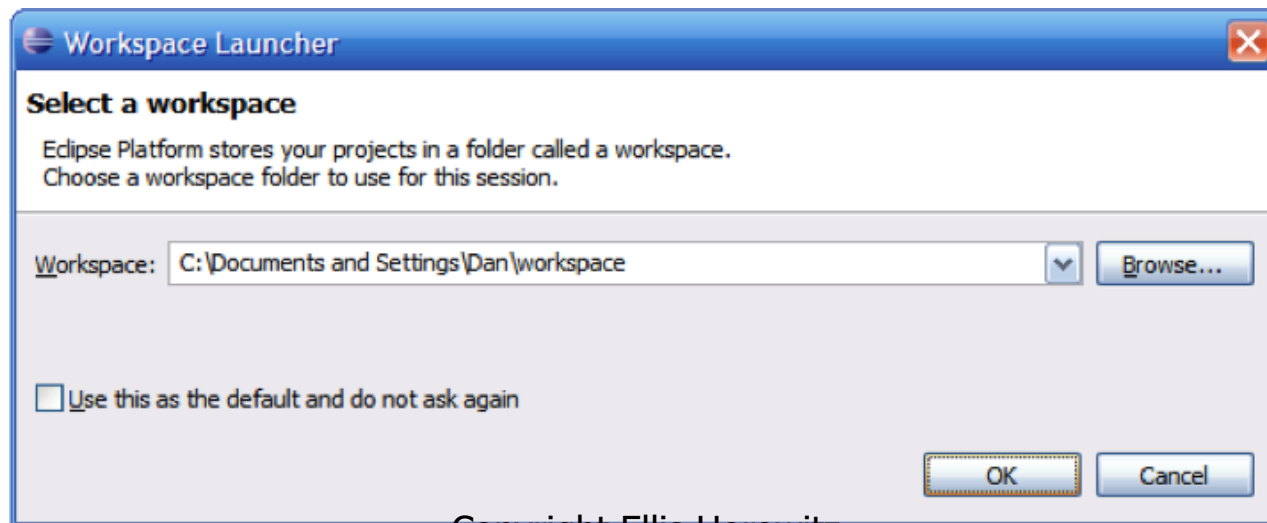
Launching Eclipse

- Once you have the environment setup, go ahead and launch eclipse
- You should see a splash screen such as the one below...

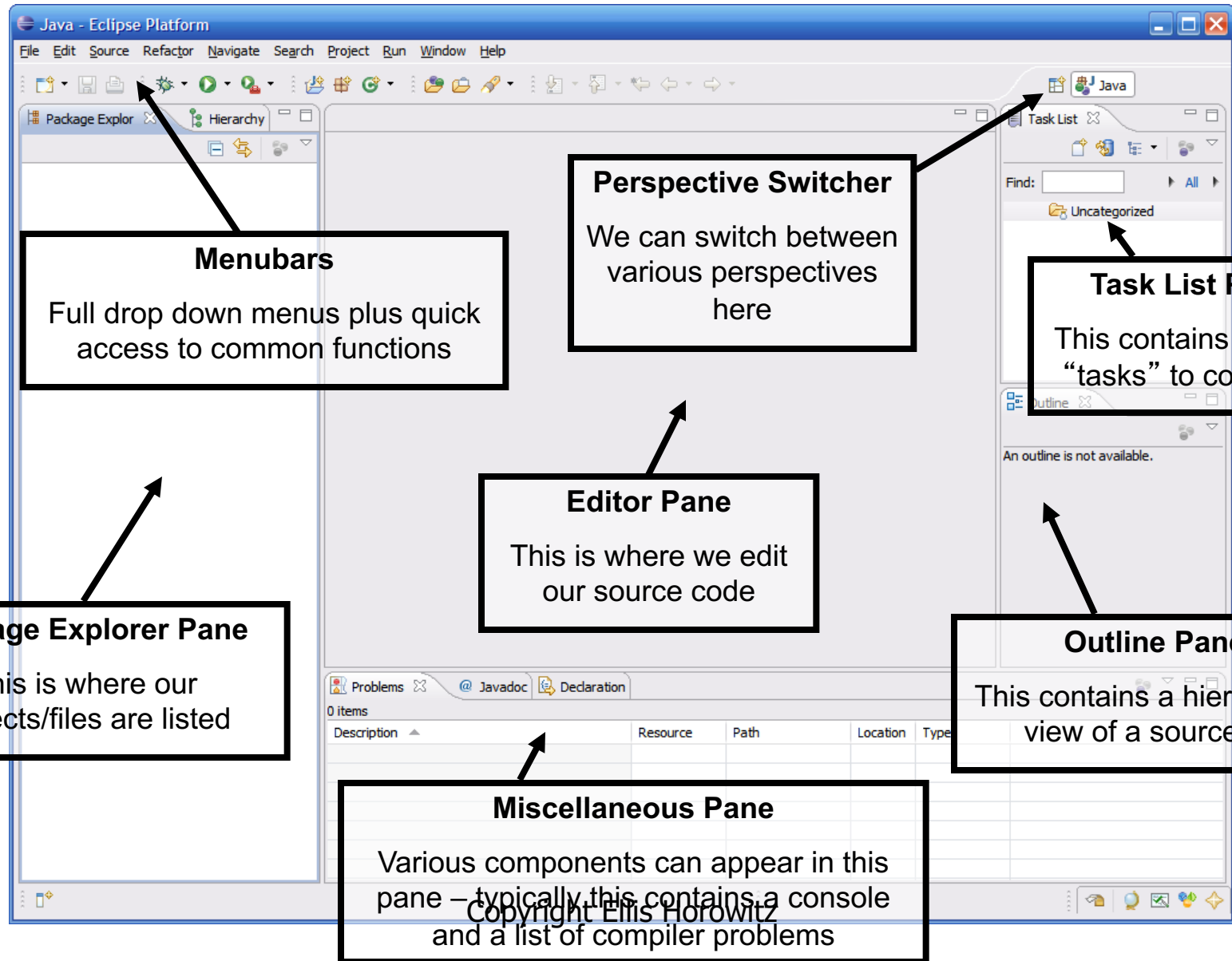


Selecting a Workspace

- In Eclipse, all of your code will live under a *workspace*
- A *workspace* is nothing more than a location where we will store the source code and where Eclipse will write out preferences
- Eclipse allows you to have multiple workspaces – each tailored in its own way
- Choose a location where you want to store your files, then click OK

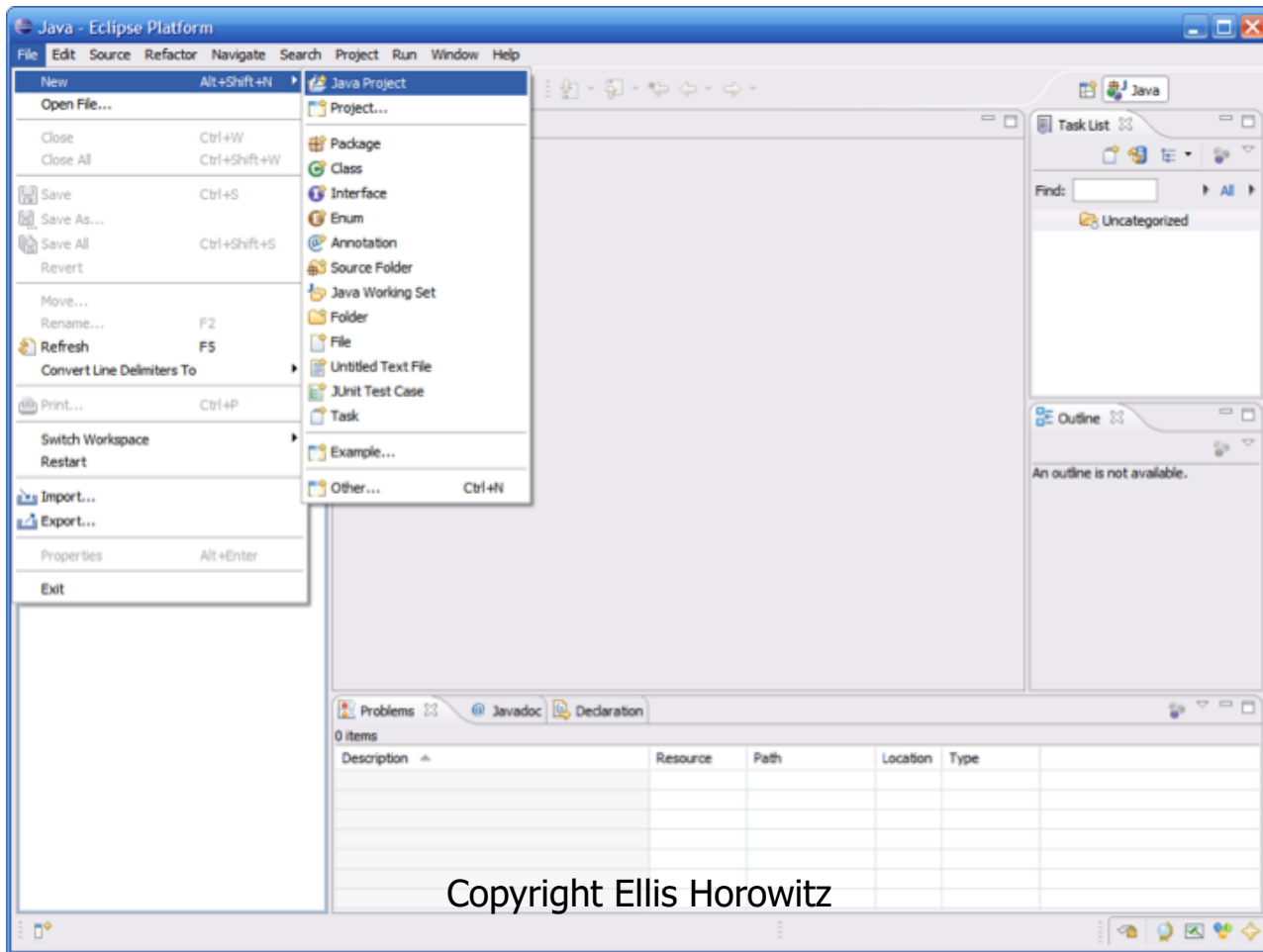


Eclipse IDE Components



Creating a New Project

- All code in Eclipse needs to live under a project
- To create a project: File → New → Java Project



Creating a New Project (continued)

- Enter a name for the project, then click Finish

New Java Project

Create a Java project in the workspace or in an external location.

Project name:

Contents

☒ Create new project in workspace
☐ Create project from existing source

Directory:

JRE

☒ Use default JRE (Currently 'jre1.6.0_05') [Configure JREs...](#)
☐ Use a project specific JRE:
☐ Use an execution environment JRE:

Project layout

☐ Use project folder as root for sources and class files
☒ Create separate folders for sources and class files [Configure default...](#)

Working sets

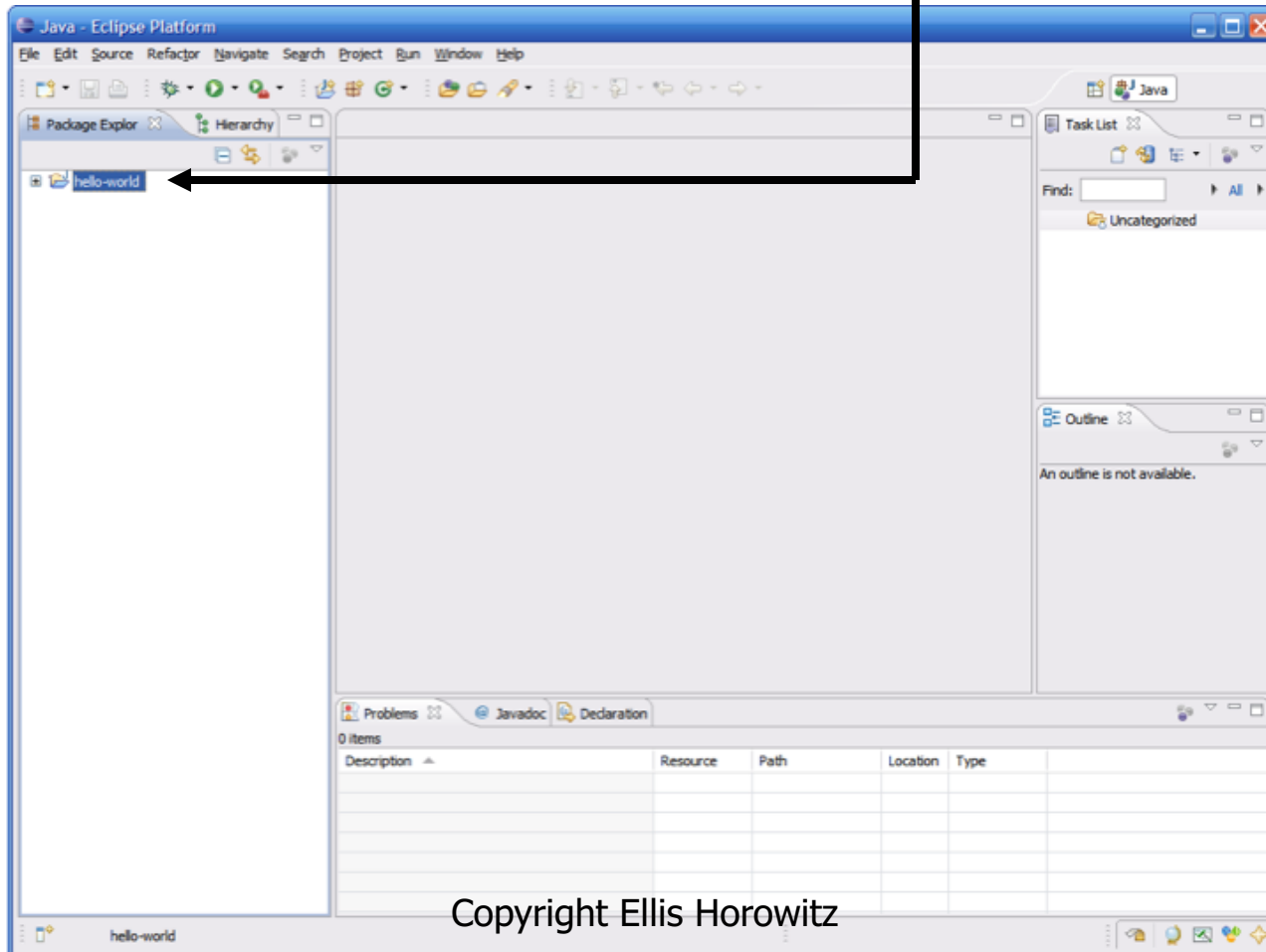
☐ Add project to working sets

Working sets:

Hello-world Project

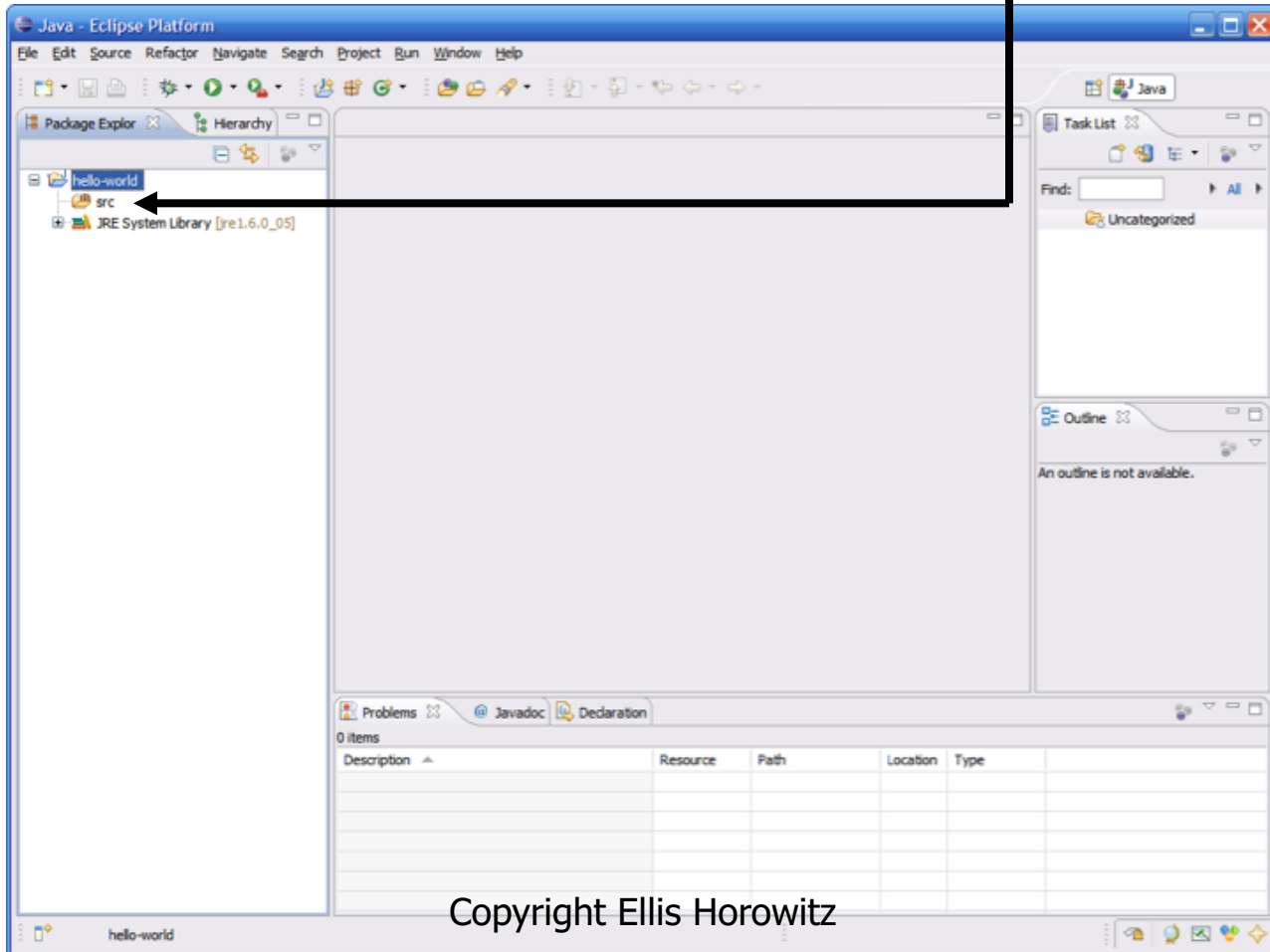
Creating a New Project (continued)

- The newly created project should then appear under the Package Explorer



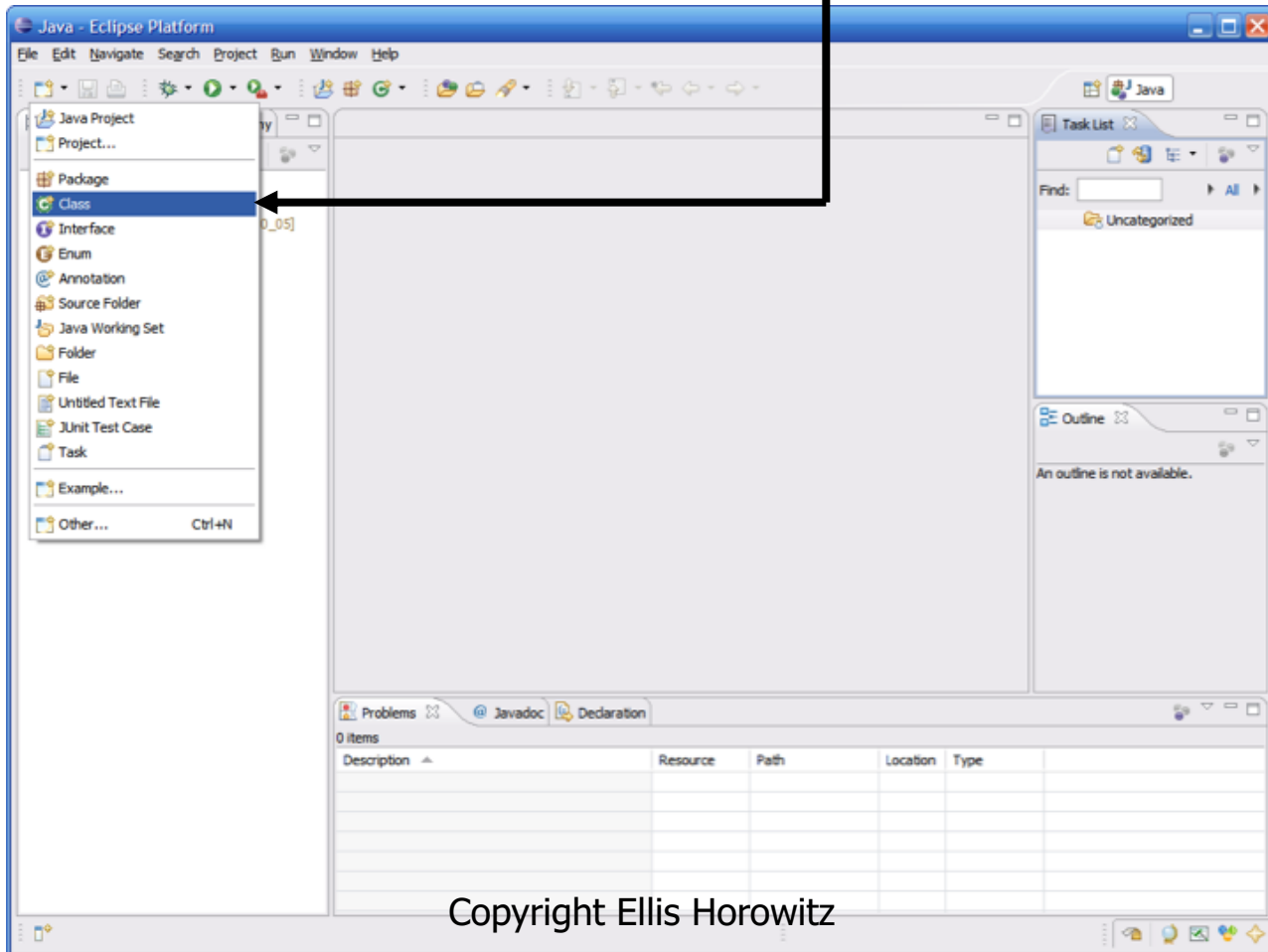
The src folder

- Eclipse automatically creates a folder to store your source code in called src

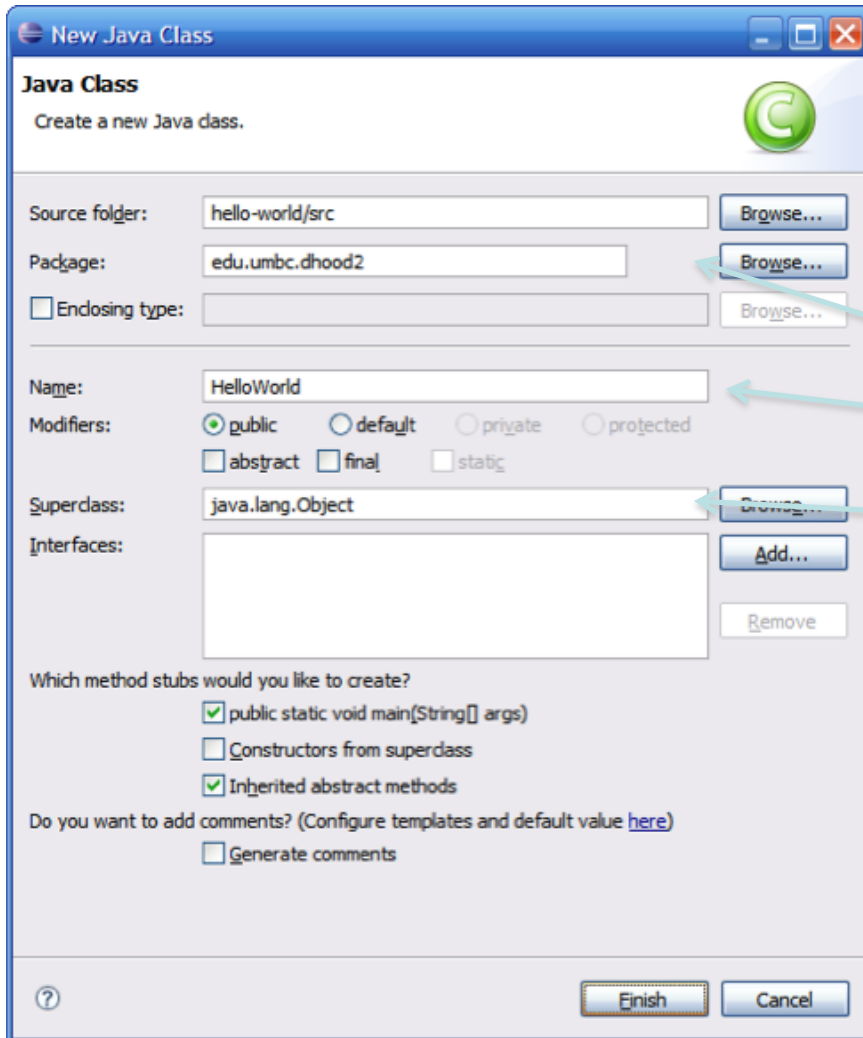


Creating a Class

- To create a class, simply click on the New button, then select Class



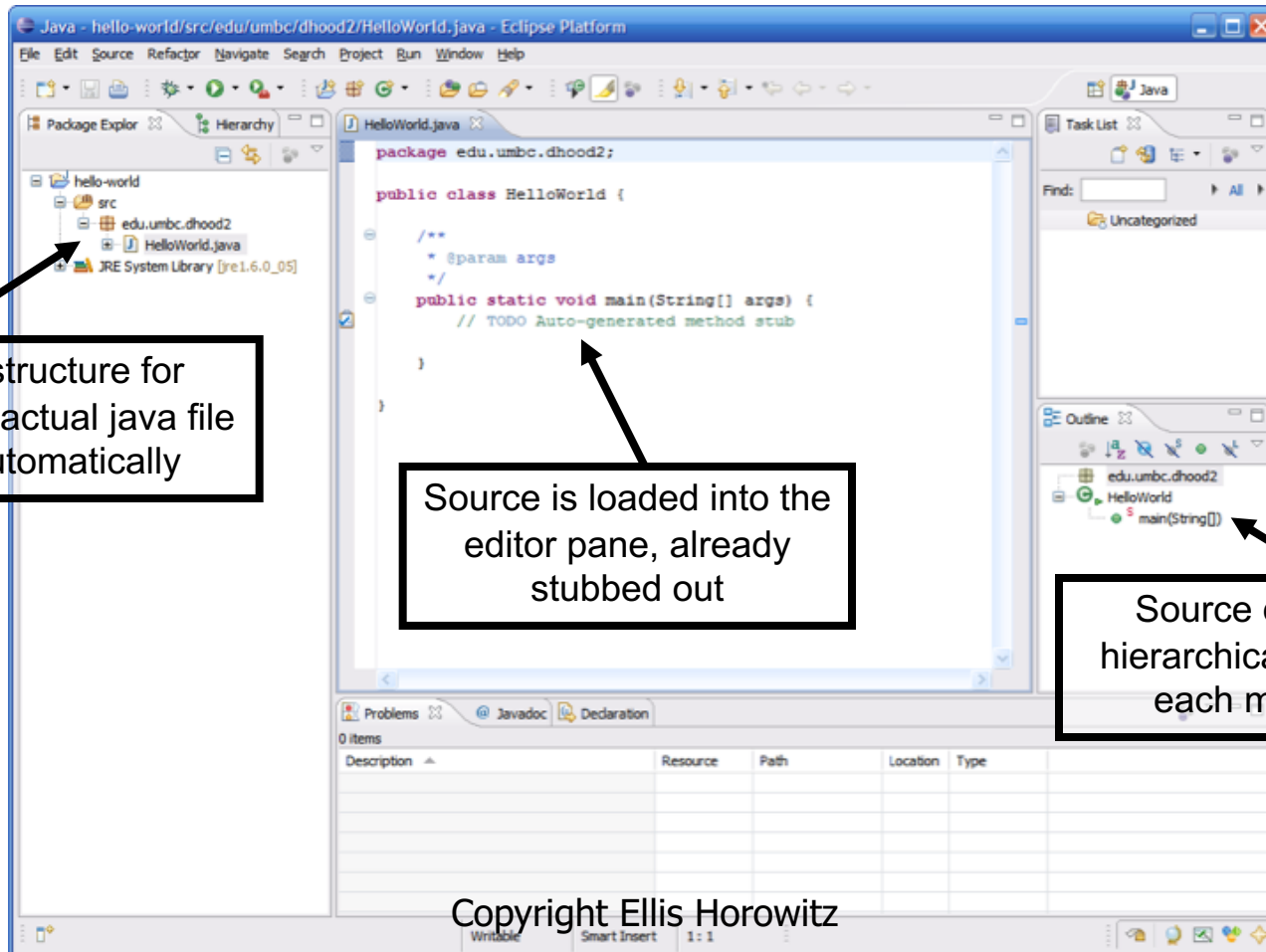
Creating a Class (continued)



- This brings up the new class wizard
- From here you can specify the following...
 - Package
 - Class name
 - Superclass
 - Whether or not to include a main
 - Etc...
- Fill in necessary information then click Finish to continue

The Created Class

- As you can see a number of things have now happened...

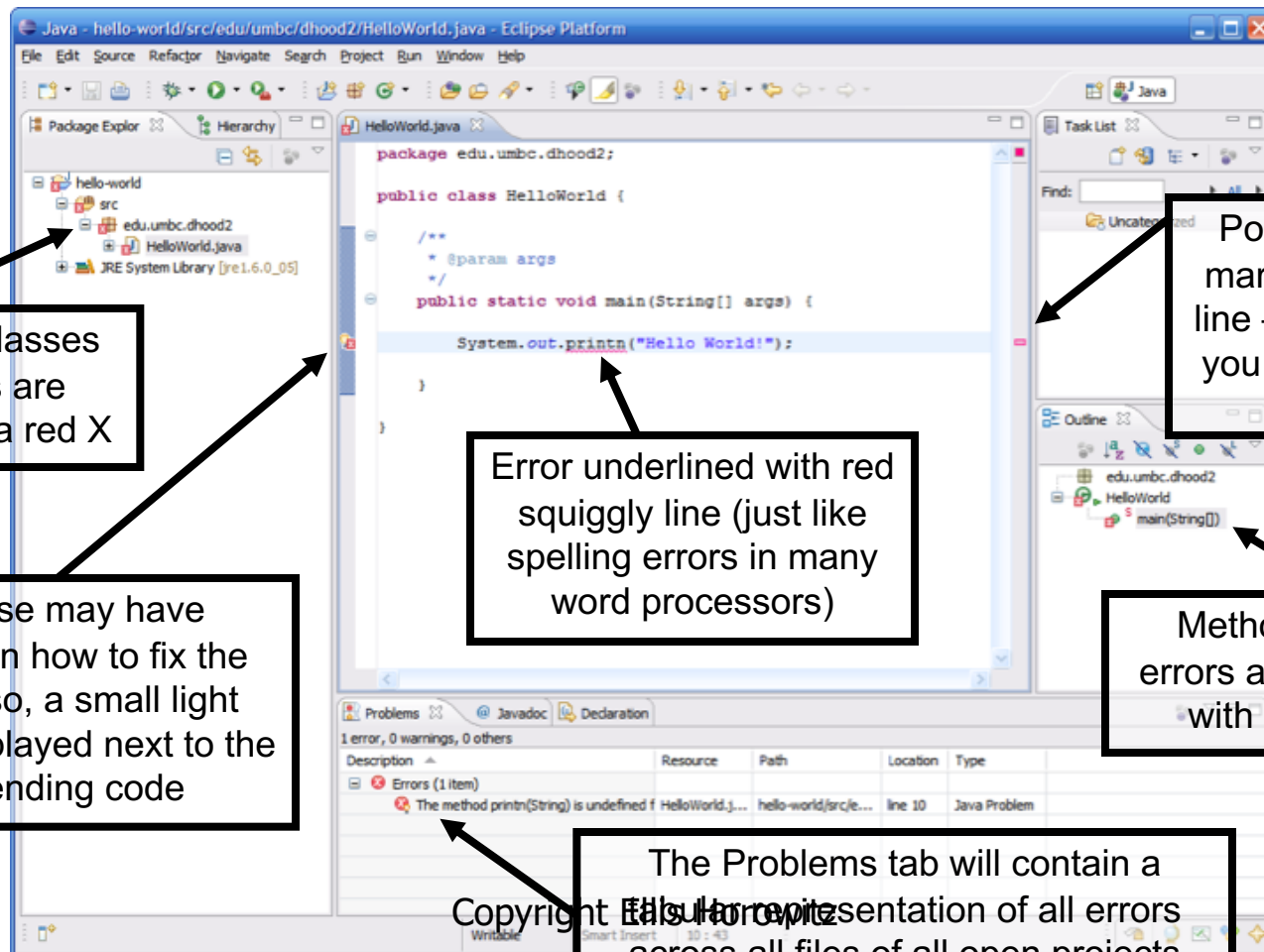


Compiling Source Code

- One important feature of Eclipse is that it automatically compiles your code in the background
- This means that errors can be corrected when made
 - We all know that iterative development is an excellent approach to developing code, but going to shell to do a compile can interrupt the normal course of development
 - You no longer need to go to the command prompt and compile code directly

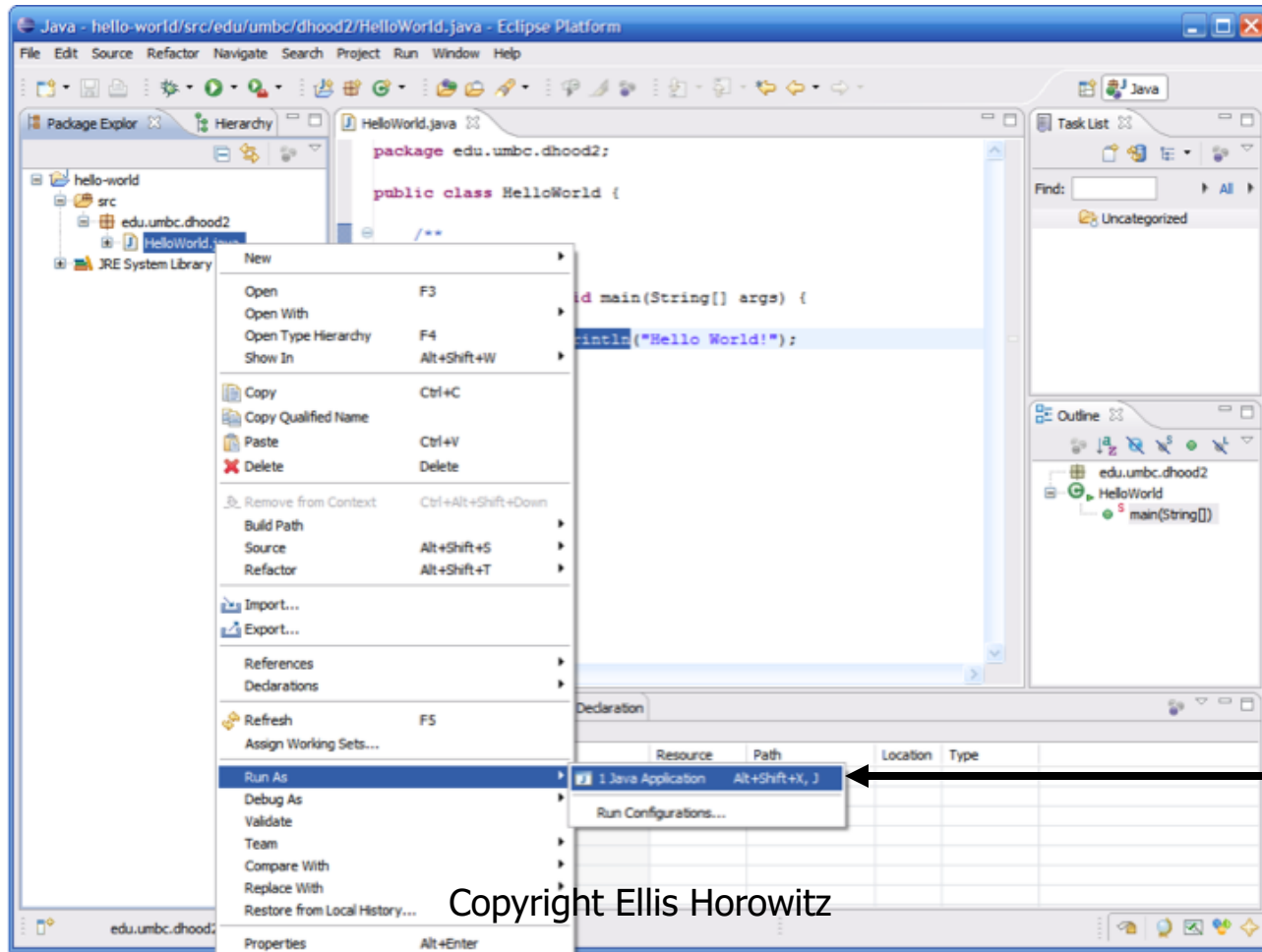
Example Compilation Error

- This code contains a typo in the println statement...



Running Code

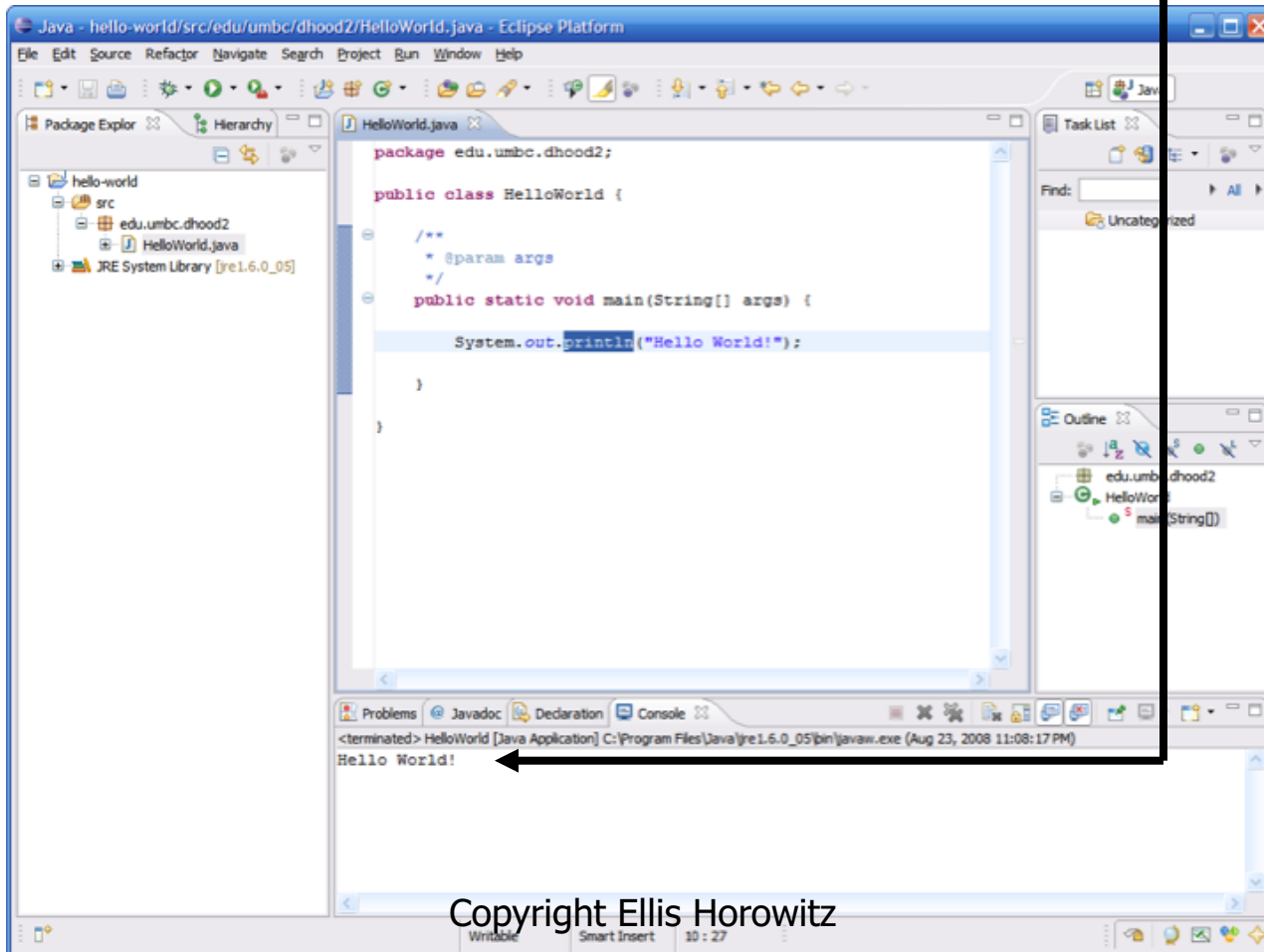
- An easy way to run code is to right click on the class and select Run As → Java Application



Copyright Ellis Horowitz

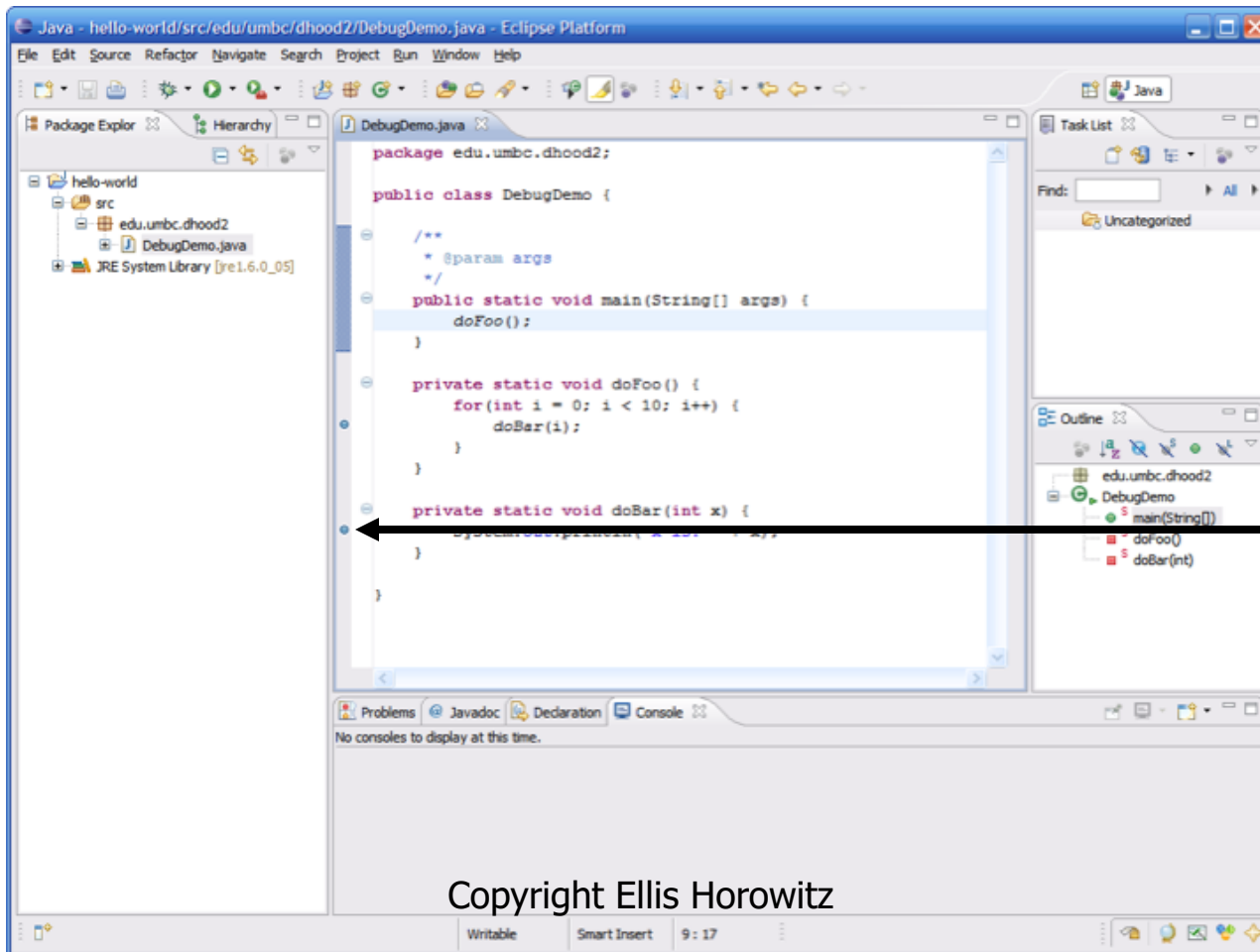
Running Code (continued)

- The output of running the code can be seen in the Console tab in the bottom pane



Debugging Code

- Eclipse comes with a pretty good built-in debugger
- You can set break points in your code by double clicking in the left hand margin – break points are represented by these blue bubbles



End of Eclipse Tutorial

Tools for Surface Web Crawling

- **Command line for issuing http requests**
 - wget, pre-installed in Ubuntu
 - get a single page
 - wget `http://www.example.com/index.html`
 - support http, ftp etc., e.g.
 - wget `ftp://ftp.gnu.org/pub/gnu/wget/wget-latest.tar.gz`
 - curl, OSX pre-installed also supports http requests
- **Simple crawling programs**
 - Crawler4j, written in Java
 - Scrapy: <http://scrapy.org>, written in Python
- **Large-scale crawling programs**
 - Heritrix, crawler for archive.org
 - Nutch, Apache Software Foundation

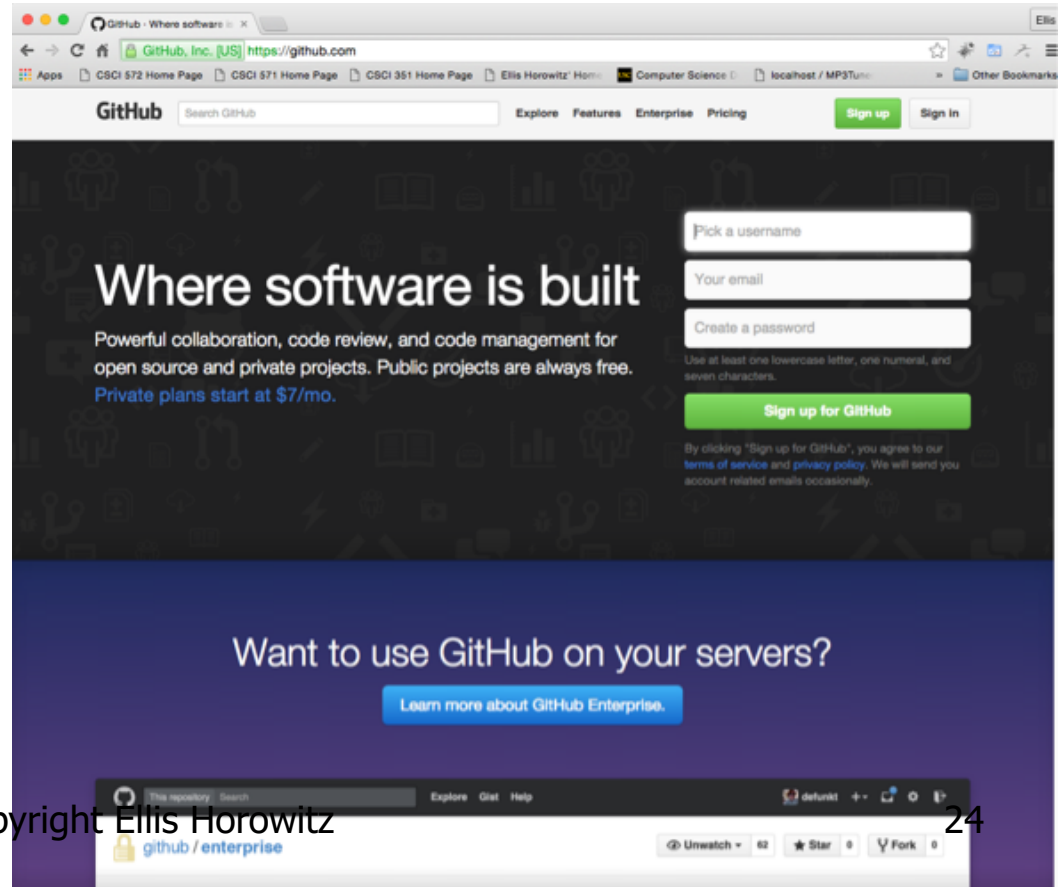
How To Get a Web Page in Java

```
import java . net . * ;
import java . io . * ;
public class URLReader {
public static void main(String [] args) throws Exception { } }
    URL oracle = new URL("http://www.oracle.com/");
    BufferedReader in = new BufferedReader (
new InputStreamReader(oracle.openStream())) ;
    String inputLine ;
    while (( inputLine = in . readLine ()) != null)
        System . out . println ( inputLine );
    in . close ();
}
```

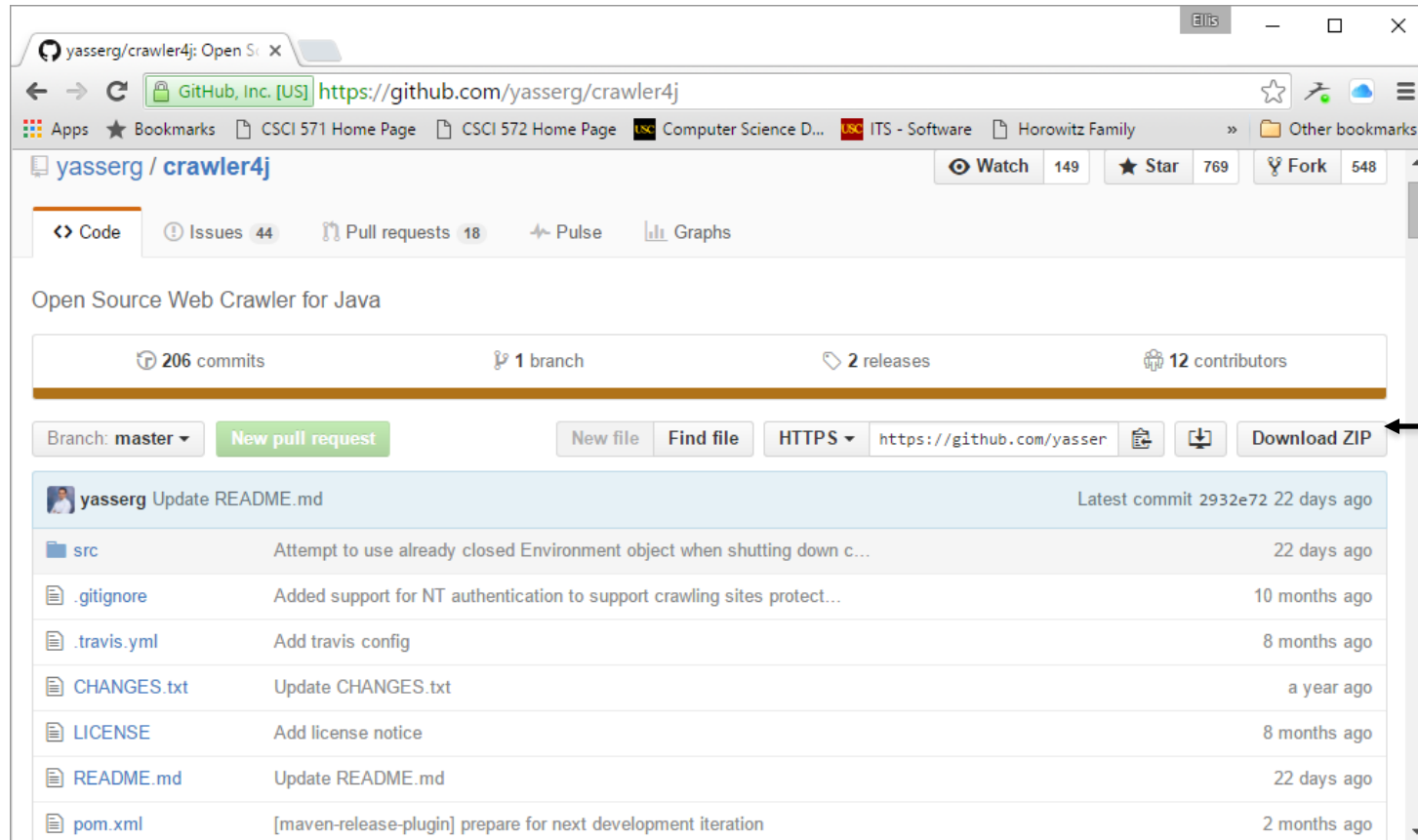
- After you create a URL, you can call the URL's `openStream()` method to get a stream from which you can read the contents of the URL.
- The `openStream()` method returns a [java.io.InputStream](#) object

Instructions for Installing Crawler4j

- download crawler4j from github
 - **GitHub** is a web-based repository hosting service for software. Originally the Git system offered distributed revision control and source code management (SCM) functionality, but on the command line; GitHub offers a web interface and some additional features.
 - As of Dec 2016, GitHub reports having over 24 million users and over 35 million repositories



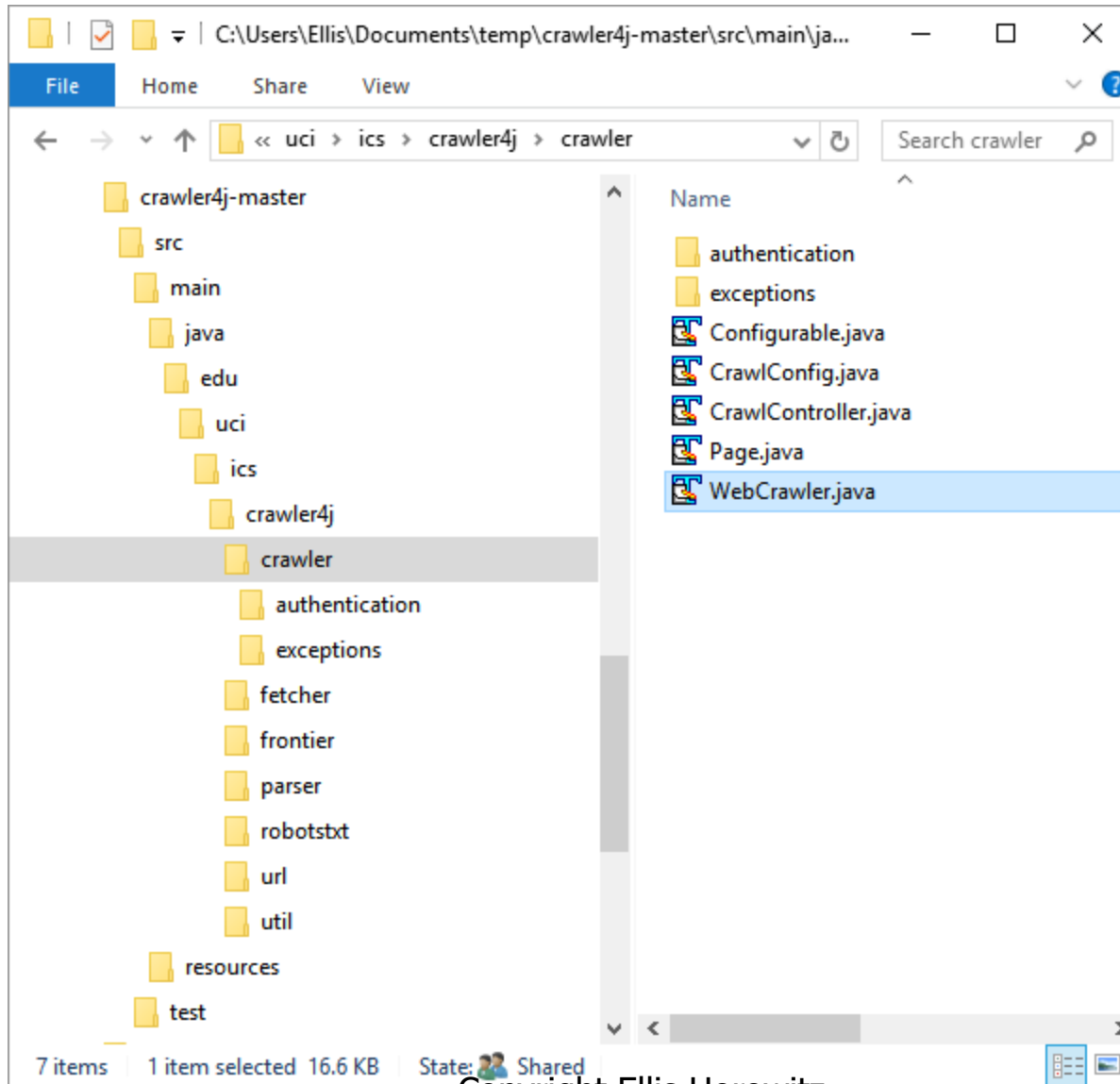
Downloading Crawler4j from GitHub



See especially the README file page at
<https://github.com/yasserg/crawler4j/blob/master/README.md>

Copyright Ellis Horowitz

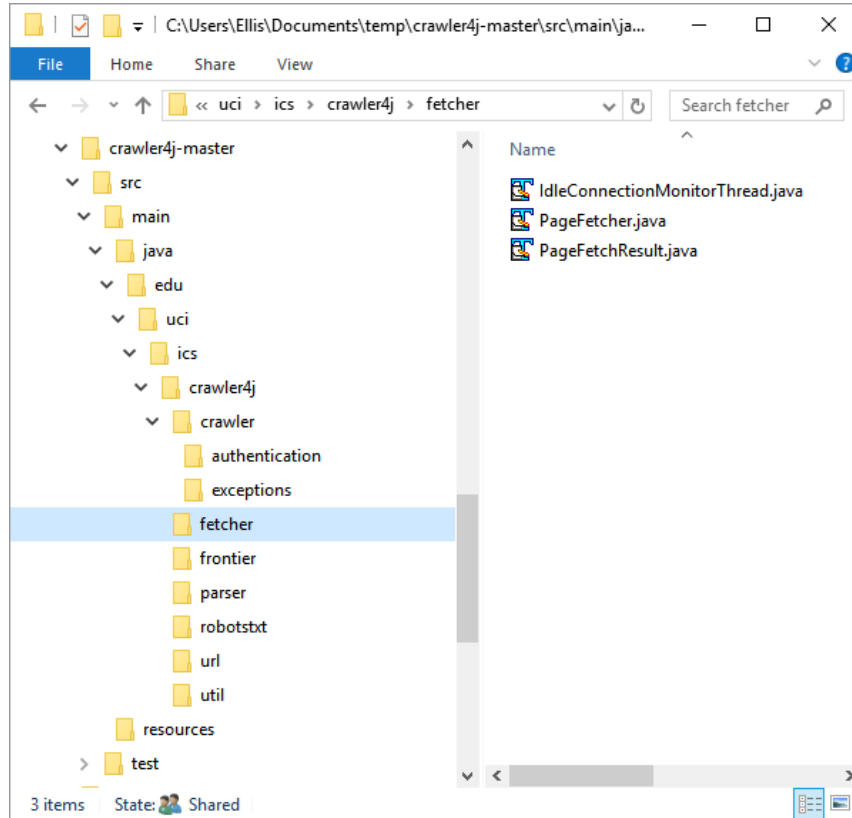
Crawler4j Source Code



Copyright Ellis Horowitz

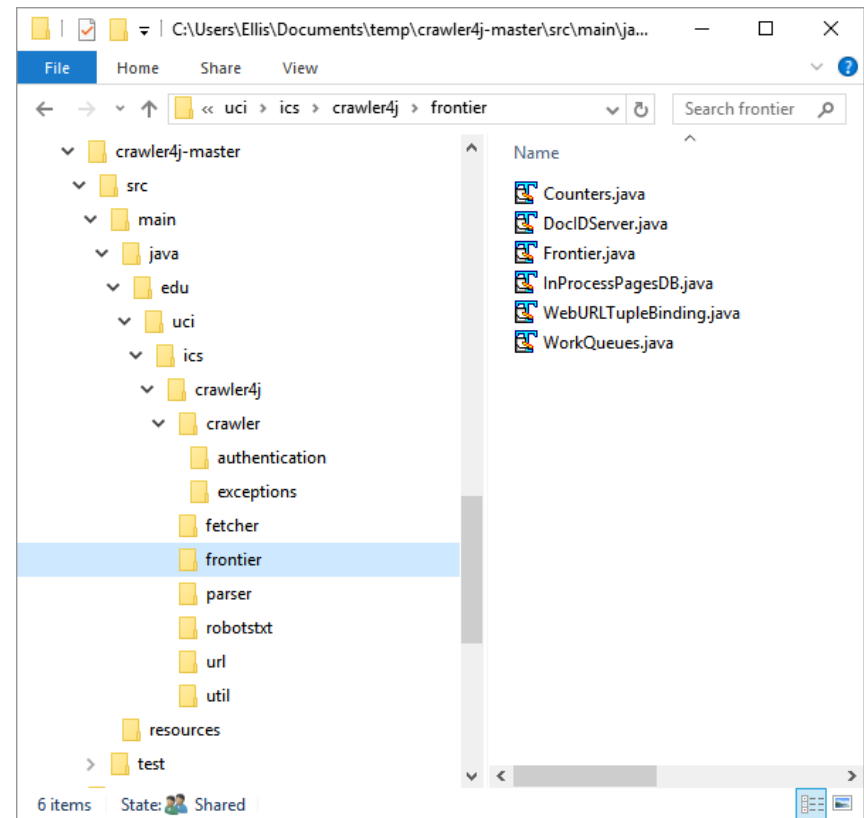
Crawler folder, a good place to start; look especially at `WebCrawler.java`

Crawler4j Source code is Logically Organized into folders



Fetcher Code handles:

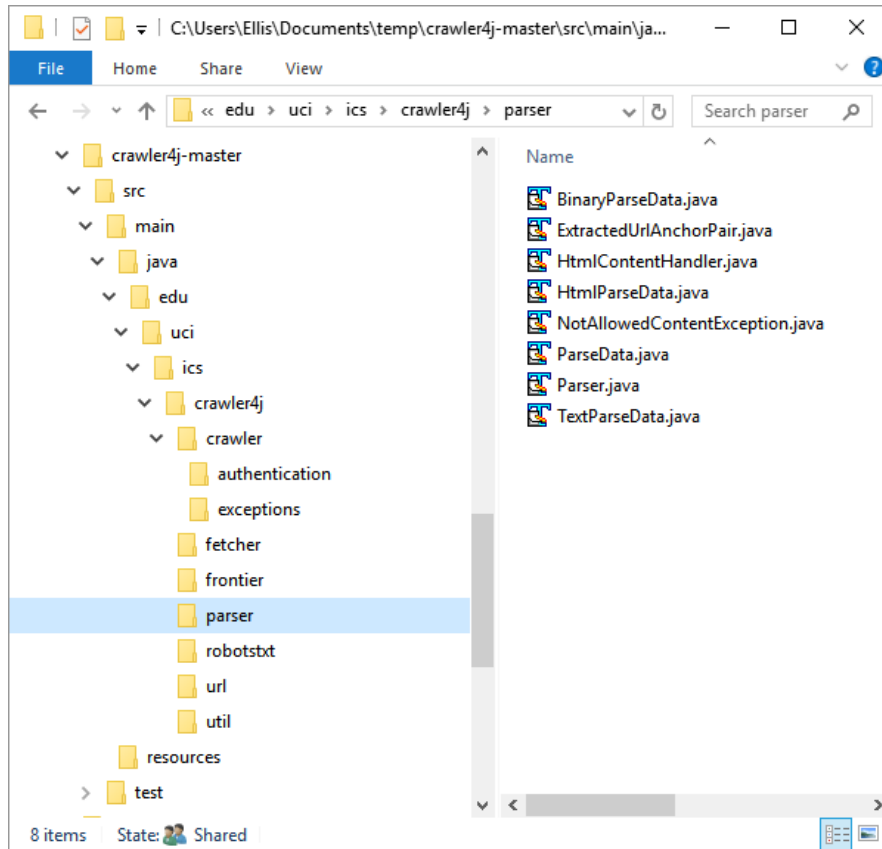
- schemes: http, https
- politeness delay;
- redirects;
- max-size settings;
- expired connections



Frontier Code handles:

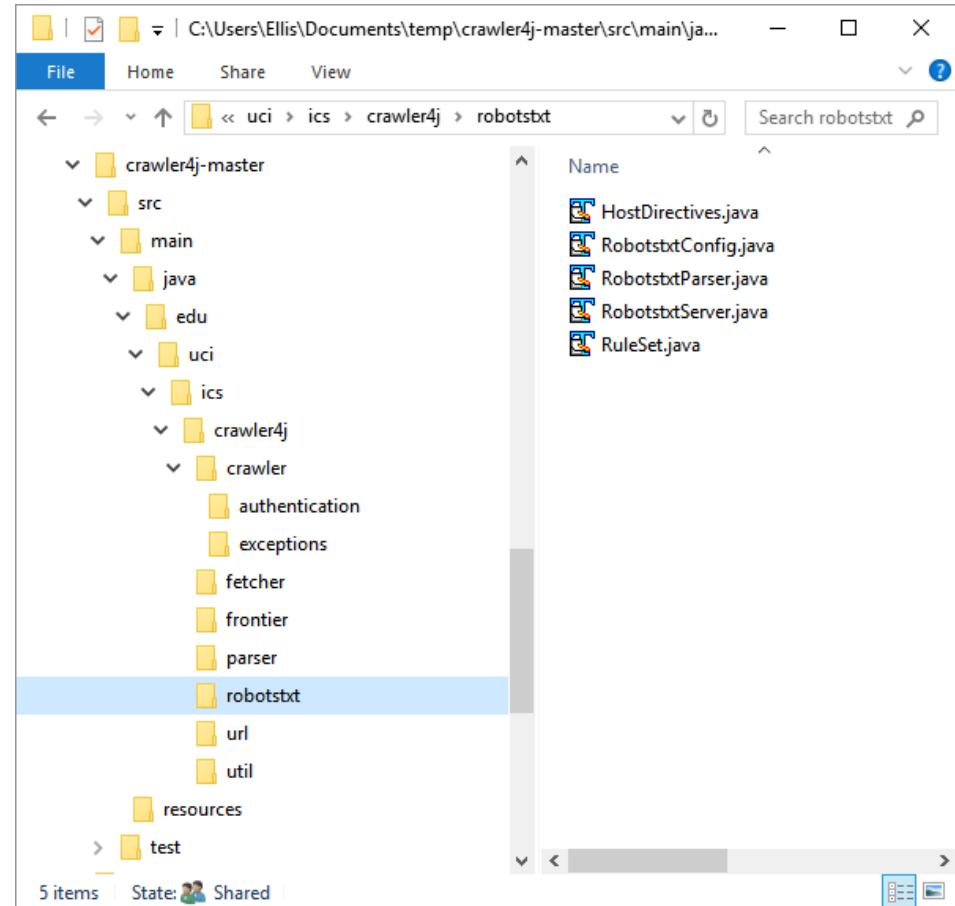
- statistics database;
- previously seen URLs
- queue of pending URLs

Crawler4j Routines are Named According to their Function



Parser Code handles:

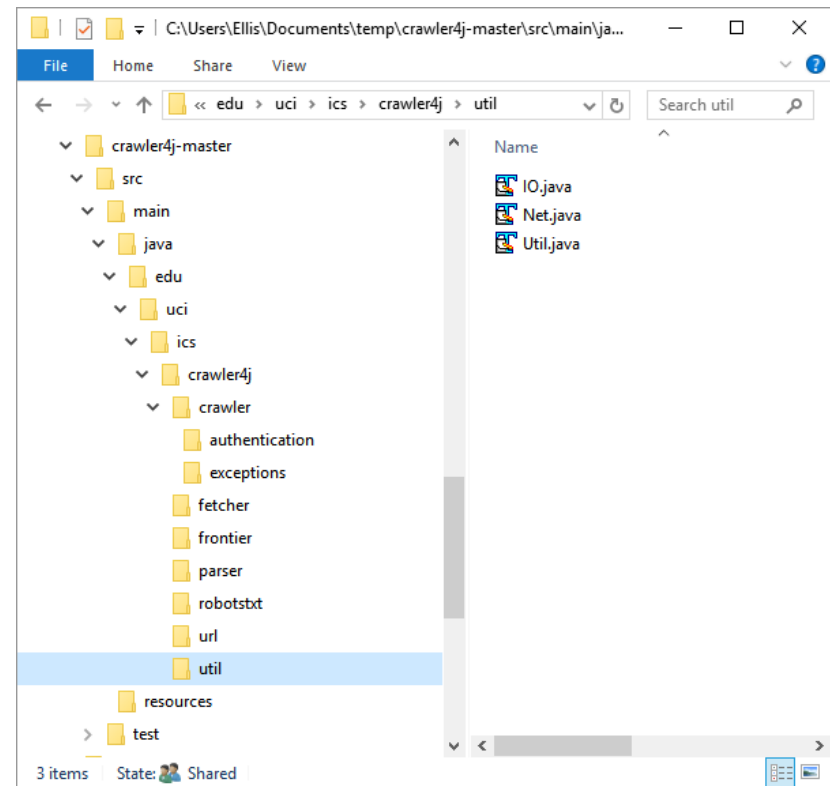
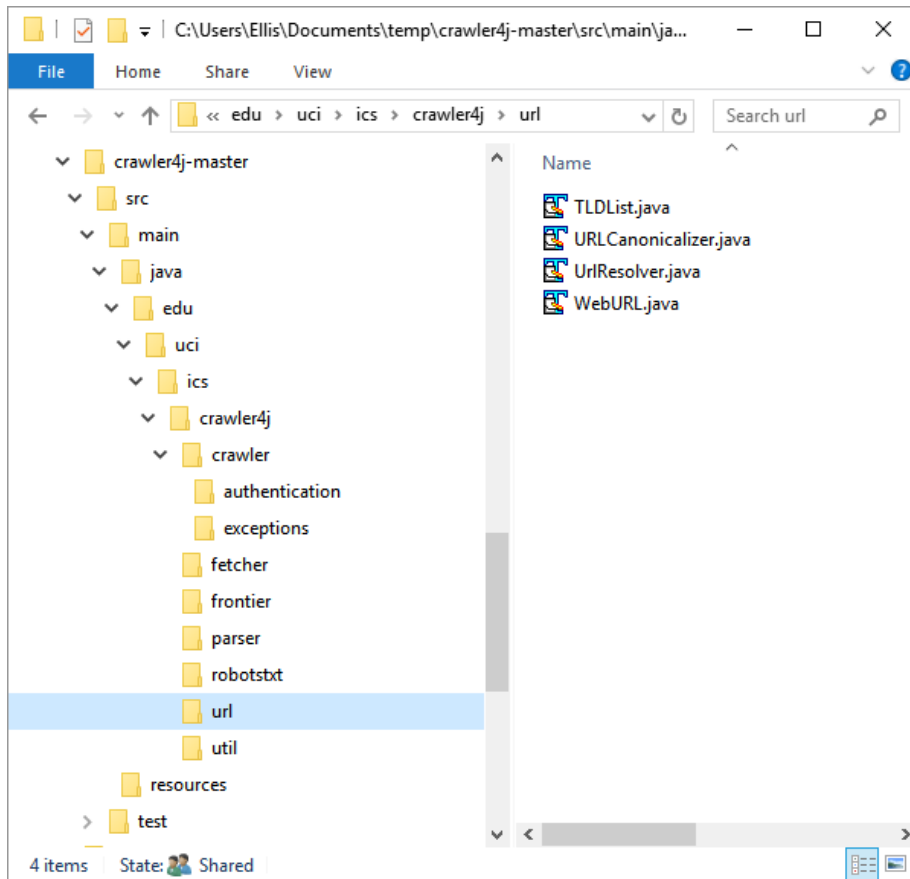
- binary data
- html pages
- extracting links



Robots.txt Code handles:

- fetching and re-fetching robots.txt
- caching robots.txt files
- interpreting commands
- working with Page Fetcher

More crawler4j Source code

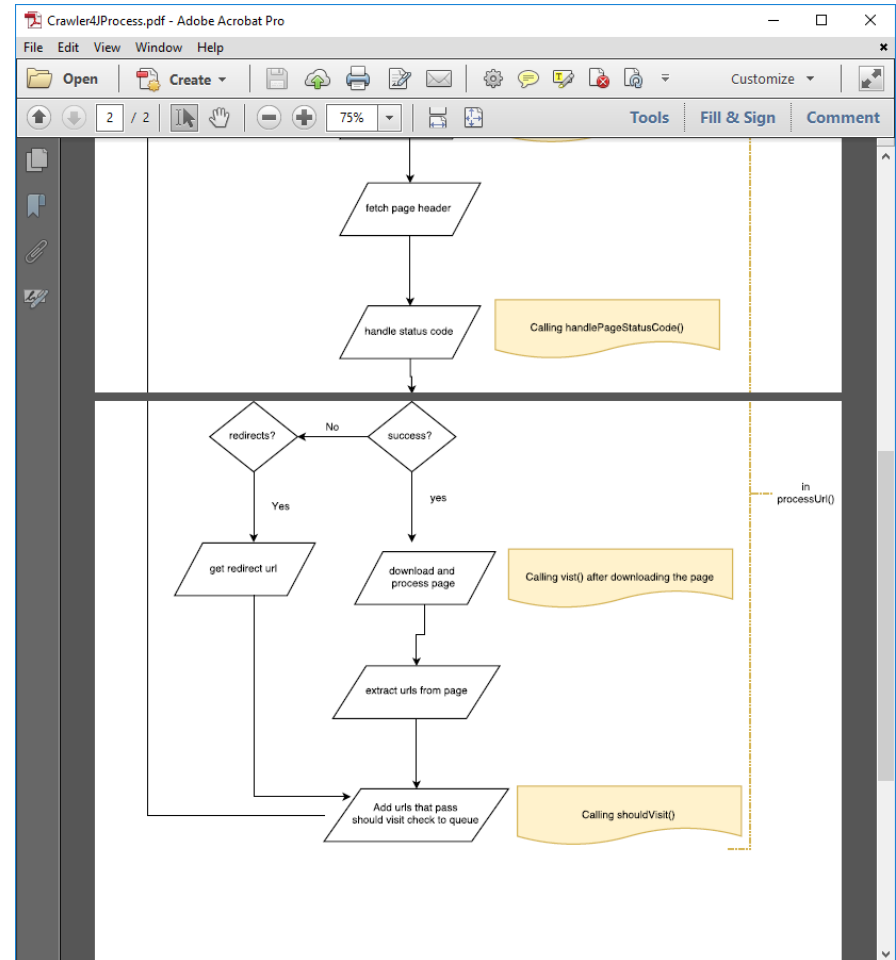
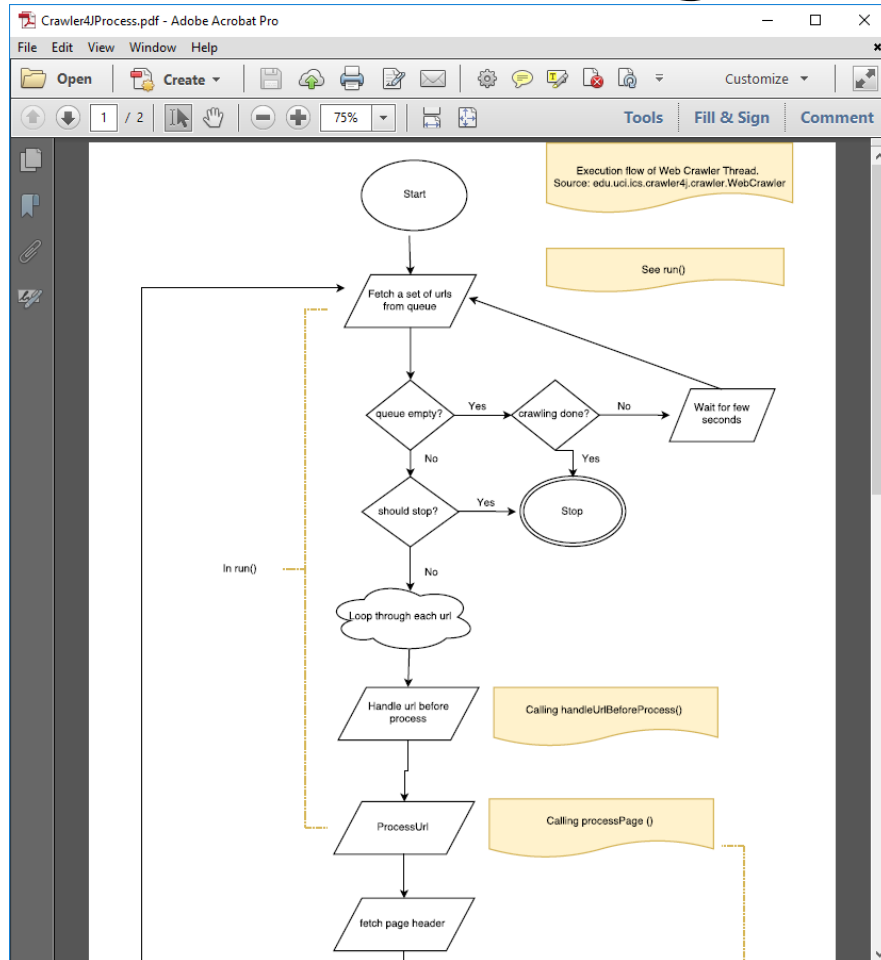


utility routines

URL resolver and canonicalizer handles:

- checking against list of TLDs
- normalizes URL, removes . or .., etc
- alters name/value pairs
- converts #nn values
- evaluates <base>

Logic Flowchart



<http://www-scf.usc.edu/~csci572/2018Spring/hw2/Crawler4JProcess.pdf>

Configuring the Crawler and Seeding it

```
public class Controller {  
    public static void main(String[] args) throws Exception {  
        String crawlStorageFolder = "/data/crawl";  
        int numberOfCrawlers = 7;  
        CrawlConfig config = new CrawlConfig();  
        config.setCrawlStorageFolder(crawlStorageFolder);  
        /* Instantiate the controller for this crawl.*/  
        PageFetcher pageFetcher = new PageFetcher(config);  
        RobotstxtConfig robotstxtConfig = new RobotstxtConfig();  
        RobotstxtServer robotstxtServer = new RobotstxtServer(robotstxtConfig, pageFetcher);  
        CrawlController controller = new CrawlController(config, pageFetcher, robotstxtServer);  
        /* For each crawl, you need to add some seed urls. These are the first  
        * URLs that are fetched and then the crawler starts following links  
        * which are found in these pages */  
        controller.addSeed("http://www.cnn.com/");  
        /* Start the crawl. This is a blocking operation, meaning that your code  
        * will reach the line after this only when crawling is finished. */  
        controller.start(MyCrawler.class, numberOfCrawlers);  
    }  
}
```

folder to store
downloads;

#crawlers

set up pagefetcher
and robots.txt
handler

crawling
www.cnn.com

Defining Which Pages to Crawl

```
public class MyCrawler extends WebCrawler {  
    private final static Pattern FILTERS =  
Pattern.compile(".*(\\.(css|js|gif|jpg" + "|png|mp3|mp3|zip|gz))$");  
    /** This method receives two parameters. The first parameter is the page  
    * in which we have discovered this new url and the second parameter is  
    * the new url. You should implement this function to specify whether  
    * the given url should be crawled or not (based on your crawling logic).  
    * In this example, we are instructing the crawler to ignore urls that  
    * have css, js, git, ... extensions and to only accept urls that start  
    * with "http://www.cnn.com/". In this case, we didn't need the  
    * referring Page parameter to make the decision. */  
    @Override  
    public boolean shouldVisit(Page referringPage, WebURL url) {  
        String href = url.getURL().toLowerCase();  
        return !FILTERS.matcher(href).matches()  
            && href.startsWith("http://www.cnn.com/");  
    }  
}
```

see next slide

Matching URLs

- `".*(\\.(css|js|gif|jpg" + "|png|mp3|mp4|zip|gz))$"`
- A regular expression, specified as a string, must first be compiled into an instance of this class.
- a Matcher object that can match arbitrary character sequences against the regular expression
- See <https://docs.oracle.com/javase/7/docs/api/java/util/regex/Pattern.html>
- In the above there are two strings concatenated by plus; consider the simpler form:
- `".*(\\.(css|js|zip|gz))$"`
- `.` matches any character
- `*` matches zero or more of preceding character
- `\\.` matches a literal dot
- `$` anchors the pattern at the end of the string

Parsing the Downloaded Page

```
/** This function is called when a page is fetched and ready
 * to be processed by your program. */
@Override
public void visit(Page page) {
    String url = page.getWebURL().getURL();
    System.out.println("URL: " + url);
    if (page.getParseData() instanceof HtmlParseData) {
        HtmlParseData htmlParseData = (HtmlParseData) page.getParseData();
        String text = htmlParseData.getText();
        String html = htmlParseData.getHtml();
        Set<WebURL> links = htmlParseData.getOutgoingUrls();
        System.out.println("Text length: " + text.length());
        System.out.println("Html length: " + html.length());
        System.out.println("Number of outgoing links: " + links.size());
    }
}
```

The Actual Exercise

- *the URLs it attempts to fetch, **fetch.csv***. The number of rows should be no more than 20,000 as that is our pre-set limit.
- *the files it successfully downloads, **visit.csv***; clearly the number of rows will be less than the number of rows in *fetch.csv*
- *all of the URLs that were discovered and processed in some way; **urls.csv***. This file could be much larger than 20,000 rows as it will have numerous repeated URLs

Things to Save

- Fetch statistics:
 - # fetches attempted:
The total number of URLs that the crawler attempted to fetch. This is usually equal to the MAXPAGES setting if the crawler reached that limit; less if the website is smaller than that.
 - # fetches succeeded:
The number of URLs that were successfully downloaded in their entirety, i.e. returning a HTTP status code of 2XX.
 - # fetches failed or aborted:
The number of fetches that failed for whatever reason, including, but not limited to: HTTP redirections (3XX), client errors (4XX), server errors (5XX) and other network-related errors.
-

Outgoing URLs

- Outgoing URLs: statistics about URLs extracted from visited HTML pages
 - Total URLs extracted:
The grand total number of URLs extracted from all visited pages
 - # unique URLs extracted:
The number of *unique* URLs encountered by the crawler
 - # unique URLs within the news web site:
The number of *unique* URLs encountered that are associated with the news website,
i.e. the URL begins with the given root URL of the news website.
 - # unique URLs outside the news website:
The number of *unique* URLs encountered that were *not* from the website.

Sample Crawl Report for News Day Using 20,000 as the Download Limit

News site crawled: <https://www.newsday.com/>

Fetch Statistics

=====

fetches attempted: 19998
fetches succeeded: 15370
fetches aborted: 0
fetches failed: 4628

Outgoing URLs

=====

Total URLs extracted: 424225
unique URLs extracted: 93445
unique URLs within News Site: 31535
unique URLs outside News Site: 61910

Status Codes

=====

200 OK: 15370
301 Moved Permanently: 4435
401 Unauthorized: 0
403 Forbidden: 0
404 Not Found: 158

File Sizes

=====

< 1KB: 23
1KB ~ <10KB: 555
10KB ~ <100KB: 13676
100KB ~ <1MB: 1116
>= 1MB: 0

Content Types

=====

text/html: 13343
image/gif: 1
image/tif: 0
image/jpeg: 1886
image/png: 21
application/pdf: 0

Sample Fetch File for News Day

fetch_NewsDay - Excel															
File Home Insert Page Layout Formulas Data Review View Tell me what you want to do...															
Clipboard Font Alignment Number Conditional Formatting Styles Cells															
A1	https://www.newsday.com/														
A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	
1	https://www.newsday.com/	200													
2	https://www.newsday.com/lifestyle/family/learn-to-paint-parties-	200													
3	https://www.newsday.com/long-island/politics	200													
4	https://www.newsday.com/entertainment/theater	200													
5	https://www.newsday.com/business/federal-reserve-imposes-new	200													
6	https://www.newsday.com/classifieds/real-estate/best-places-to-l	200													
7	https://www.newsday.com/towns	301													
8	https://www.newsday.com/long-island/education	200													
9	https://www.newsday.com/homedelivery/terms	302													
10	https://www.newsday.com/long-island/politics/spin-cycle	200													
11	https://www.newsday.com/sports/basketball/nets/lakers-nets-1.1	200													
12	https://www.newsday.com/onlineaccess	302													
13	https://www.newsday.com/entertainment/movies/winchester-revi	200													
14	https://www.newsday.com/long-island/long-island-and-new-york-c	200													
15	https://www.newsday.com/long-island/i-love-ny-signs-pulled-1.16	200													
16	https://www.newsday.com/sports/football/giants	200													
17	https://www.newsday.com/services/profile-7.386?modify=1	200													
18	https://www.newsday.com/classifieds/cars	200													
19	https://www.newsday.com/sports/baseball/mets	200													
20	https://www.newsday.com/business/glen-mill-apartments-robert-r	200													
21	https://www.newsday.com/long-island/suffolk/epcal-luminati-sola	200													
22	https://www.newsday.com/build/newsday-front.min.css?v=n1811	200													
23	https://www.newsday.com/business/stock-quotes-1.5276093	200													
24	https://www.newsday.com/eedition	301													
25	https://www.newsday.com/sports/football/giants/eli-manning-gia	200													
26	https://www.newsday.com/long-island/transportation/lirr-amtrak-	200													
27	https://www.newsday.com/news/nation	200													
28	https://www.newsday.com/travel/nyc-weekend-picks-our-best-be	200													
29	https://www.newsday.com/img/newsday/apple-touch-icon-precor	200													
fetch_NewsDay															

Sample Visit File for News Day

visit_NewsDay - Excel

File Home Insert Page Layout Formulas Data Review View Tell me what you want to do...

Clipboard Font Alignment Number Styles

Normal Bad Neutral Calculati

A1 https://www.newsday.com/

	A	B	C	D	E	F	G	H	I	J
1	https://www.newsday.com/	118.81640625 kb	323	text/html						
2	https://www.newsday.com/lifestyle/family/learn-to-paint-	77.1015625 kb	194	text/html						
3	https://www.newsday.com/long-island/politics	50.4560546875 kb	172	text/html						
4	https://www.newsday.com/entertainment/theater	51.31640625 kb	167	text/html						
5	https://www.newsday.com/business/federal-reserve-impos	73.822265625 kb	175	text/html						
6	https://www.newsday.com/classifieds/real-estate/best-pla	84.8583984375 kb	183	text/html						
7	https://www.newsday.com/long-island/education	50.6103515625 kb	176	text/html						
8	https://www.newsday.com/long-island/politics/spin-cycle	48.9248046875 kb	170	text/html						
9	https://www.newsday.com/sports/basketball/nets/lakers-n	83.6015625 kb	181	text/html						
10	https://www.newsday.com/entertainment/movies/winches	68.2626953125 kb	163	text/html						
11	https://www.newsday.com/long-island/long-island-and-nev	212.41015625 kb	323	text/html						
12	https://www.newsday.com/long-island/i-love-ny-signs-pulle	74.79296875 kb	178	text/html						
13	https://www.newsday.com/sports/football/giants	85.6484375 kb	239	text/html						
14	https://www.newsday.com/services/profile-7.386?modify=	33.6552734375 kb	111	text/html						
15	https://www.newsday.com/classifieds/cars	76.8466796875 kb	226	text/html						
16	https://www.newsday.com/sports/baseball/mets	77.3671875 kb	241	text/html						
17	https://www.newsday.com/business/glen-mill-apartments-	64.8115234375 kb	167	text/html						
18	https://www.newsday.com/long-island/suffolk/epcal-lumin	77.529296875 kb	178	text/html						
19	https://www.newsday.com/build/newsday-front.min.css?v=	64.123046875 kb	0	text/css						
20	https://www.newsday.com/business/stock-quotes-1.52760	60.8974609375 kb	168	text/html						
21	https://www.newsday.com/sports/football/giants/eli-manr	84.857421875 kb	181	text/html						
22	https://www.newsday.com/long-island/transportation/lirr-	77.0947265625 kb	182	text/html						
23	https://www.newsday.com/news/nation	96.8212890625 kb	279	text/html						
24	https://www.newsday.com/travel/nyc-weekend-picks-our-	91.3759765625 kb	197	text/html						
25	https://www.newsday.com/img/newsday/apple-touch-icon	1.7080078125 kb	0	image/png						
26	https://www.newsday.com/sports/college/hofstra/boogie-	90.3564453125 kb	186	text/html						
27	https://www.newsday.com/sports/football/super-bowl/sup	79.814453125 kb	185	text/html						
28	https://www.newsday.com/entertainment/movies	50.7431640625 kb	174	text/html						
29	https://www.newsday.com/classifieds/cars/consumer-repo	94.8076171875 kb	203	text/html						

visit_NewsDay

What to Submit

- Compress all of the above into a single zip archive and name it:
crawl.zip
- Use only standard zip format. Do **NOT** use other formats such as zipx, rar, ace, etc. For example the zip file might contain the following three files:
 1. CrawlReport_Newsday.txt,
 2. fetch_Newsday.csv
 3. visit_Newsday.csv
- To submit your file electronically to the csci572 account enter the following command from your UNIX prompt:
- `$ submit -user csci572 -tag hw2 crawl.zip`